

assignment_06_BasakAtanu.R

atanu

2022-05-20

```
# Assignment: ASSIGNMENT 6
# Name: Basak, Atanu
# Date: 2022-05-02

## Set the working directory to the root of your DSC 520 directory

setwd("C:\\Users\\atanu\\Documents\\BellevueUniversity_MSDS\\DSC520\\Repository\\dsc520")

## Load the `data/r4ds/heights.csv` to
heights_df <- read.csv("data\\r4ds\\heights.csv")

## Load the ggplot2 library
library(ggplot2)

## Fit a linear model using the `age` variable as the predictor and `earn` as the outcome
head(heights_df)

##      earn  height    sex ed age  race
## 1 50000 74.42444  male 16  45 white
## 2 60000 65.53754 female 16  58 white
## 3 30000 63.62920 female 16  29 white
## 4 50000 63.10856 female 16  91 other
## 5 51000 63.40248 female 17  39 white
## 6  9000 64.39951 female 15  26 white

age_lm <- lm(earn~age, data=heights_df)

## View the summary of your model using `summary()`
summary(heights_df)

##      earn      height      sex      ed
## Min.   : 200    Min.   :57.50  Length:1192  Min.   : 3.0
## 1st Qu.:10000   1st Qu.:64.01   Class :character  1st Qu.:12.0
## Median :20000   Median :66.45   Mode  :character  Median :13.0
## Mean   :23155   Mean   :66.92                Mean   :13.5
## 3rd Qu.:30000   3rd Qu.:69.85                3rd Qu.:16.0
## Max.   :200000   Max.   :77.05                Max.   :18.0
##      age      race
## Min.   :18.00  Length:1192
## 1st Qu.:29.00  Class :character
## Median :38.00  Mode  :character
```

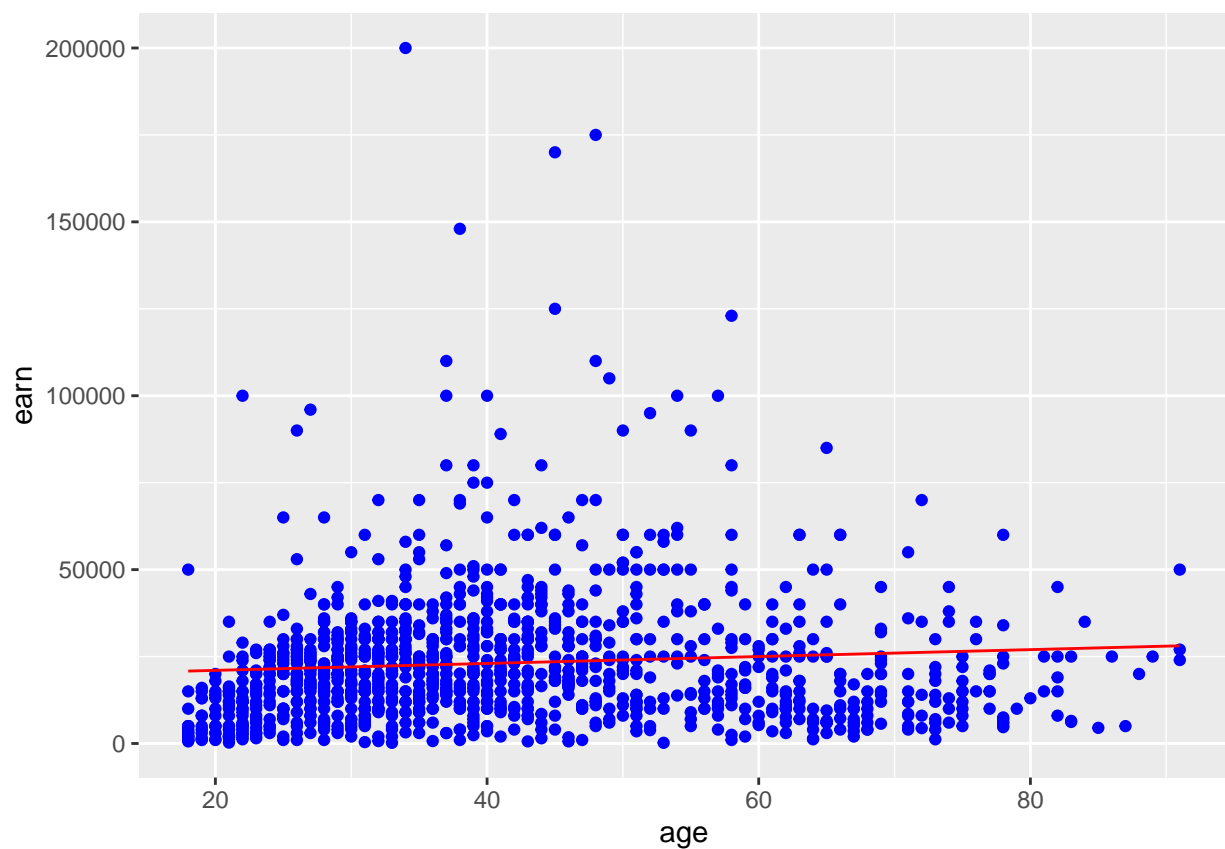
```
## Mean    :41.38
## 3rd Qu. :51.00
## Max.    :91.00
```

```
## Creating predictions using `predict()`
```

```
age_predict_df <- data.frame(earn = predict(age_lm, data.frame(age=heights_df$age)), age=heights_df$age)
#head(age_predict_df)
```

```
## Plot the predictions against the original data
```

```
ggplot(data = heights_df, aes(y = earn, x = age)) +
  geom_point(color='blue') +
  geom_line(color='red', data = age_predict_df, aes(y=earn, x=age))
```



```
mean_earn <- mean(heights_df$earn)
## Corrected Sum of Squares Total
sst <- sum((mean_earn - heights_df$earn)^2)
## Corrected Sum of Squares for Model
ssm <- sum((mean_earn - age_predict_df$earn)^2)
## Residuals
residuals <- heights_df$earn - age_predict_df$earn
#residuals
## Sum of Squares for Error
sse <- sum(residuals^2)
```

```
## R Squared  $R^2 = SSM/SST$ 
summary(age_lm)
```

```
##
## Call:
## lm(formula = earn ~ age, data = heights_df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -25098 -12622  -3667   6883 177579
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) 19041.53    1571.26   12.119 < 2e-16 ***
## age          99.41       35.46    2.804 0.00514 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 19420 on 1190 degrees of freedom
## Multiple R-squared:  0.006561,    Adjusted R-squared:  0.005727
## F-statistic:  7.86 on 1 and 1190 DF,  p-value: 0.005137
```

```
r_squared <- summary(age_lm)$r.squared
r_squared
```

```
## [1] 0.006561482
```

```
## Number of observations
n <- nrow(heights_df)
## Number of regression parameters
p <- 2
## Corrected Degrees of Freedom for Model (p-1)
dfm <- p-1
## Degrees of Freedom for Error (n-p)
dfe <- n-p
## Corrected Degrees of Freedom Total:  DFT = n - 1
dft <- n-1

## Mean of Squares for Model:  MSM = SSM / DFM
msm <- ssm/dfm
msm
```

```
## [1] 2963111900
```

```
## Mean of Squares for Error:  MSE = SSE / DFE
mse <- sse / dfe
mse
```

```
## [1] 376998968
```

```
## Mean of Squares Total:   $MST = SST / DFT$ 
mst <- sst / dft
mst
```

```
## [1] 379170348
```

```
## F Statistic  $F = MSM/MSE$ 
f_score <- msm/mse
f_score
```

```
## [1] 7.859735
```

```
## Adjusted R Squared  $R^2 = 1 - (1 - R^2)(n - 1) / (n - p)$ 
adjusted_r_squared <- 1 - (1 - r_squared)*(n - 1) / (n - p)
adjusted_r_squared
```

```
## [1] 0.005726659
```

```
## Calculate the p-value from the F distribution
p_value <- pf(f_score, dfm, dft, lower.tail=F)
```